

# Machine Precision

- ▶ Computer numbers (floating point numbers) are a **finite** subset of rational numbers.
- ▶ There is a smallest positive computer number  $\epsilon$  so that

$$1 + \epsilon > 1$$

# Machine Precision

```
>> eps
ans =
2.220446049250313e-16

>> x = 1 + eps; disp([eps, (x-1)/eps])
2.220446049250313e-16    1.000000000000000e+00

>> delta = 0.75*eps
delta =
1.665334536937735e-16

>> x = 1 + delta; disp([delta, (x - 1)/delta])
1.665334536937735e-16    1.333333333333333e+00

>> Delta = 0.5*eps
Delta =
1.110223024625157e-16

>> x=1+Delta;disp([Delta, (x - 1)/Delta])
1.110223024625157e-16          0
```

# Overflow

```
>> x=2^1023
x=2^1023

x =
8.9885e+307

>> x = 2*x
x = 2*x

x =
Inf

>> y = 2^(-1023)
y = 2^(-1023)

y =
1.1125e-308

>> y = y/(2^51)
y = y/(2^51)

y =
4.9407e-324

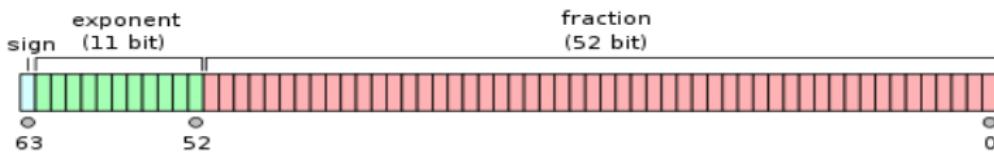
>> y = y / 2
y = y / 2

y =
0
```

# Apple Memory Chip: Huge but finite capacity



# IEEE 754 Double Precision Format



Every floating point number fits within 64 bits.

# Round-off Errors and Floating Point Arithmetic

- ▶ **Binary Machine Numbers:** any double precision non-zero *floating point number* has form

$$x = (-1)^s 2^{c-1023} (1 + f), \quad \text{with 64 bits.}$$

- ▶  $s$  = SIGN BIT: 0 for  $x > 0$  and 1 for  $x < 0$ .
- ▶  $c$  = CHARACTERISTIC, with 11 bits:

$$c = c_1 \cdot 2^{10} + c_2 \cdot 2^9 + c_3 \cdot 2^8 + c_4 \cdot 2^7 + c_5 \cdot 2^6 + c_6 \cdot 2^5 + c_7 \cdot 2^4 + c_8 \cdot 2^3 + c_9 \cdot 2^2 + c_{10} \cdot 2^1 + c_{11} \cdot 2^0,$$

- with each  $c_j = 0$  or 1.
- ▶  $f$  = MANTISSA with 52 bits

$$f = f_1 \cdot \left(\frac{1}{2}\right) + \cdots + f_{52} \cdot \left(\frac{1}{2}\right)^{52} = \sum_{j=1}^{52} f_j \cdot \left(\frac{1}{2}\right)^j, \quad \text{each } f_j = 0 \text{ or } 1.$$

- ▶ **Floating Point:** Binary point always comes after 1, independent of  $c$ .

# Round-off Errors and Floating Point Arithmetic

- **Binary Machine Numbers:** Example binary string

0 10000000011 10111001000100

- $s = 0, c = (10000000011)_2 = 1024 + 2 + 1 = 1027$ , and

$$f = 1 \cdot \left(\frac{1}{2}\right)^1 + 1 \cdot \left(\frac{1}{2}\right)^3 + 1 \cdot \left(\frac{1}{2}\right)^4 + 1 \cdot \left(\frac{1}{2}\right)^5 + 1 \cdot \left(\frac{1}{2}\right)^8 + 1 \cdot \left(\frac{1}{2}\right)^{12}.$$

$$(-1)^s 2^{e-1023} (1+f) = (-1)^0 \cdot 2^{1027-1023} \left( 1 + \left( \frac{1}{2} + \frac{1}{8} + \frac{1}{16} + \frac{1}{32} + \frac{1}{256} + \frac{1}{4096} \right) \right)$$

$$= 27.56640625.$$

# Round-off Errors and Floating Point Arithmetic

- ▶  **$k$ -digit Decimal Machine Numbers:**

$$x = \pm 0.d_1 d_2 \cdots d_k \times 10^n, \quad \text{where} \quad 1 \leq d_1 \leq 9, \quad 0 \leq d_i \leq k, i \geq 2.$$

- ▶ Any positive real number

$$\begin{aligned}y &= 0.d_1 d_2 \cdots d_k d_{k+1} d_{k+2} \cdots \times 10^n, \\&\approx 0.d_1 d_2 \cdots d_k \times 10^n \stackrel{\text{def}}{=} fl(y) \quad (\textbf{chopping}) \\&\approx 0.\delta_1 \delta_2 \cdots \delta_k \times 10^n \stackrel{\text{def}}{=} fl(y) \quad (\textbf{rounding}),\end{aligned}$$

where

$$\textbf{rounding} = \textbf{chopping on } y + 5 \times 10^{n-(k+1)}.$$

- ▶ If  $d_{k+1} < 5$ : **rounding = chopping**.
- ▶ If  $d_{k+1} \geq 5$ : cut off  $d_{k+1}$  and below, then add 1 to  $d_k$ .

# Round-off Errors and Floating Point Arithmetic

- ▶ 5-digit Decimal Machine Numbers for  $\pi$ :

$$\pi = 0.314159265 \dots \times 10^1$$

$$\approx 0.31415 \times 10^1 = 3.1415 \quad (\text{chopping})$$

$$\approx (0.31415 + 0.00001) \times 10^1 = 3.1416 \quad (\text{rounding}).$$

## Absolute error vs. relative error

Suppose that  $p^*$  is an approximation to  $p \neq 0$ .

- ▶ **absolute error** =  $|p - p^*|$ ,
- ▶ **relative error** =  $\frac{|p - p^*|}{|p|}$ .

Example

- ▶ **absolute errors:**

$$|\pi - 3.1415| \approx 9 \times 10^{-5}, \quad |\pi - 3.1416| \approx 7 \times 10^{-6}.$$

- ▶ **relative errors:**

$$\frac{|\pi - 3.1415|}{\pi} \approx 3 \times 10^{-5}, \quad \frac{|\pi - 3.1416|}{\pi} \approx 2 \times 10^{-6}.$$

A cool \$200,000 wager by LeSean McCoy,  $\frac{\text{wager}}{\text{salary}} \approx 3\%$



David Payne Purdum   
 @DavidPurdum



Follow

LeSean McCoy bet \$200,000 on Warriors to win Finals. It's largest bet Planet Hollywood took on Finals. Would pay \$62,500.  
h/t [@RomanEdmond1](#)

8:54 AM - Jun 4, 2017

## Relative error for chopping

Suppose that  $y = 0.d_1d_2 \dots d_k d_{k+1} d_{k+2} \dots \times 10^n$ , with  $d_1 \geq 1$ .

$$\begin{aligned}\left| \frac{y - fl(y)}{y} \right| &= \left| \frac{0.d_1d_2 \dots d_k d_{k+1} \dots \times 10^n - 0.d_1d_2 \dots d_k \times 10^n}{0.d_1d_2 \dots \times 10^n} \right| \\ &= \left| \frac{0.d_{k+1}d_{k+2} \dots \times 10^{n-k}}{0.d_1d_2 \dots \times 10^n} \right| = \left| \frac{0.d_{k+1}d_{k+2} \dots}{0.d_1d_2 \dots} \right| \times 10^{-k}.\end{aligned}$$

But  $0.d_1d_2 \dots d_k d_{k+1} d_{k+2} \dots \geq 0.1$ ,

$$\left| \frac{y - fl(y)}{y} \right| \leq \frac{1}{0.1} \times 10^{-k} = 10^{-k+1}.$$

## Relative error for rounding

Suppose that  $y = 0.d_1d_2 \cdots d_kd_{k+1}d_{k+2} \cdots \times 10^n$ , with  $d_1 \geq 1$ .

$$\left| \frac{y - fl(y)}{y} \right| \leq 0.5 \times 10^{-k+1}.$$

**Proof:** Exercise in text.

Floating Point Arithmetic:

RELATIVE ERROR  $\approx 10^{-k+1}$  INDEPENDENT OF  $n$ .

## Machine addition, subtraction, multiplication, and division

$$x \oplus y = fl(fl(x) + fl(y)), \quad x \otimes y = fl(fl(x) \times fl(y)),$$

$$x \ominus y = fl(fl(x) - fl(y)), \quad x \oslash y = fl(fl(x) \div fl(y)).$$

**Some computations involve millions of these operations,  
the result could be very different from expected.**

## Cancellation of significant digits

Suppose that  $x$  and  $y$  do not differ much:

$$\begin{aligned}x &= 0.d_1 \cdots d_p \alpha_{p+1} \cdots \times 10^n \\&= 0.d_1 \cdots d_p \alpha_{p+1} \cdots \alpha_k \times 10^n + \epsilon_x = fl(x) + \epsilon_x, \\y &= 0.d_1 \cdots d_p \beta_{p+1} \cdots \times 10^n \\&= 0.d_1 \cdots d_p \beta_{p+1} \cdots \beta_k \times 10^n + \epsilon_y = fl(y) + \epsilon_y,\end{aligned}$$

with  $\epsilon_x, \epsilon_y \approx 10^{n-k}$ .

The floating-point form of  $x - y$  is

$$fl(fl(x) - fl(y)) = 0.\sigma_{p+1} \sigma_{p+2} \dots \sigma_k \times 10^{n-p},$$

where

$$0.\sigma_{p+1} \sigma_{p+2} \dots \sigma_k = 0.\alpha_{p+1} \alpha_{p+2} \dots \alpha_k - 0.\beta_{p+1} \beta_{p+2} \dots \beta_k.$$

Roughly, **relative error** is

$$\left| \frac{\text{error in computing } x - y}{x - y} \right| \approx \left| \frac{|\epsilon_x| + |\epsilon_y|}{fl(fl(x) - fl(y))} \right| \approx \frac{10^{n-k}}{10^{n-p}} = 10^{-(k-p)}.$$

## Quadratic formula for $ax^2 + bx + c = 0$

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \quad \text{and} \quad x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}.$$

One of  $x_1$ ,  $x_2$  faces cancellation of significant digits if

$$|4ac| \ll b^2.$$

# Solving $ax^2 + bx + c = 0$ the better way

- ▶ Compute  $\delta = \sqrt{b^2 - 4 * a * c}$
- ▶ If  $b > 0$  then

$$x_1 = \frac{-b - \delta}{2a};$$

if  $b \leq 0$  then

$$x_1 = \frac{-b + \delta}{2a}.$$

- ▶ Vieta's formula

$$x_2 = \frac{c}{a x_1}.$$

## Roots to Quadratic to Roots (I)

```
function xx = quadroot(x)
a = 1;
b = -(x(1)+x(2));
c = x(1)*x(2);
del = sqrt(b*b-4*a*c);
xx(1) = (-b+del)/(2*a);
xx(2) = (-b-del)/(2*a);
xx =xx(:);
```

## Roots to Quadratic to Roots (II)

```
>> format long e
format long e
>> x = randn(2,1)
x = randn(2,1)

x =
1.630235289164729e+00
4.888937703117894e-01

>> xx = quadroot(x)
xx = quadroot(x)

xx =
1.630235289164729e+00
4.888937703117894e-01

>> x = [randn*1e5;randn*1e-12]
x = [randn*1e5;randn*1e-12]

x =
1.034693009917860e+05
7.268851333832379e-13

>> xx = quadroot(x)
xx = quadroot(x)

xx =
1.034693009917860e+05
0
```

## Roots to Quadratic to Roots (III)

```
>> a = randn*1e-5;b = 1; c = - randn*1e-12;
a = randn*1e-5;b = 1; c = - randn*1e-12;
>> roots([a b c])
roots([a b c])

ans =

    3.295534380226372e+05
    2.938714670966580e-13

>> del = sqrt(b*b-4*a*c)
del = sqrt(b*b-4*a*c)

del =

    1

>> x(1) = (-b+del)/(2*a);x(2) = (-b-del)/(2*a)
x(1) = (-b+del)/(2*a);x(2) = (-b-del)/(2*a)

x =

    0
    3.295534380226372e+05

>> x(2) = (-b-del)/(2*a);x(1)=(c/a)/x(2)
x(2) = (-b-del)/(2*a);x(1)=(c/a)/x(2)

x =

    2.938714670966580e-13
    3.295534380226372e+05
```

# Horner's Method for Fibonacci's Problem in 1224, with his Emperor

Solve

$$f(x) = x^3 + 2x^2 + 10x - 20 = 0.$$

Fibonacci's Solution

$$x = 1 + 22 \left( \frac{1}{60} \right) + 7 \left( \frac{1}{60} \right)^2 + 42 \left( \frac{1}{60} \right)^3 + 33 \left( \frac{1}{60} \right)^4 + 4 \left( \frac{1}{60} \right)^5 + 40 \left( \frac{1}{60} \right)^6.$$

With Horner's nested sum method, let  $\tau = \left( \frac{1}{60} \right)$ :

$$x = 1 + \tau \cdot (22 + \tau \cdot (7 + \tau \cdot (42 + \tau \cdot (33 + \tau \cdot (4 + 40\tau))))).$$

## Pseudocode for Horner's Method

```
function SUM = horner(x,a)
%
% horner's method
%
n = length(a);
SUM = a(n)*ones(size(x));
for i=n-1:-1:1
    SUM = a(i) + x .* SUM;
end
return
```

## Numerical stability: a second order recursion

For any constants  $c_1$  and  $c_2$ ,

$$p_n = c_1 \left(\frac{1}{3}\right)^n + c_2 3^n,$$

is a solution to the recursive equation

$$p_n = \frac{10}{3}p_{n-1} - p_{n-2}, \quad \text{for } n = 2, 3, \dots$$

- ▶  
$$\lim_{n \rightarrow \infty} |p_n| = \begin{cases} \infty & \text{if } c_2 \neq 0, \\ 0 & \text{otherwise.} \end{cases}$$
- ▶  
$$\begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \frac{1}{8} \begin{pmatrix} 9p_0 - 3p_1 \\ 3p_1 - p_0 \end{pmatrix}, \quad \text{given } p_0, p_1.$$
- ▶ condition  $c_2 = 3p_1 - p_0 = 0$  hard to satisfy exactly in finite precision computations.

Numerical values go crazy for  $p_0 = 1, p_1 = 1/3$ .

With five-digit rounding arithmetic,

$n$	Computed $\hat{p}_n$	Correct $p_n$	Relative Error
0	$0.10000 \times 10^1$	$0.10000 \times 10^1$	
1	$0.33333 \times 10^0$	$0.33333 \times 10^0$	
2	$0.11110 \times 10^0$	$0.11111 \times 10^0$	$9 \times 10^{-5}$
3	$0.37000 \times 10^{-1}$	$0.37037 \times 10^{-1}$	$1 \times 10^{-3}$
4	$0.12230 \times 10^{-1}$	$0.12346 \times 10^{-1}$	$9 \times 10^{-3}$
5	$0.37660 \times 10^{-2}$	$0.41152 \times 10^{-2}$	$8 \times 10^{-2}$
6	$0.32300 \times 10^{-3}$	$0.13717 \times 10^{-2}$	$8 \times 10^{-1}$
7	$-0.26893 \times 10^{-2}$	$0.45725 \times 10^{-3}$	$7 \times 10^0$
8	$-0.92872 \times 10^{-2}$	$0.15242 \times 10^{-3}$	$6 \times 10^1$

## Rate of convergence: the Big O

Suppose  $\{\beta_n\}_{n=1}^{\infty}$  is a sequence known to converge to zero, and  $\{\alpha_n\}_{n=1}^{\infty}$  converges to a number  $\alpha$ . If a positive constant  $K$  exists with

$$|\alpha_n - \alpha| \leq K|\beta_n|, \quad \text{for large } n,$$

then we say that  $\{\alpha_n\}_{n=1}^{\infty}$  converges to  $\alpha$  with **rate, or order, of convergence**  $O(\beta_n)$ . (This expression is read “big oh of  $\beta_n$ ”.) It is indicated by writing  $\alpha_n = \alpha + O(\beta_n)$ . ■

# Matlab Primer

matlab primer

All Shopping Videos Images News More Settings Tools

About 261,000 results (0.32 seconds)

**MATLAB Primer - MathWorks**  
[https://www.mathworks.com/help/pdf\\_doc/matlab/getstart.pdf](https://www.mathworks.com/help/pdf_doc/matlab/getstart.pdf) ▾  
A description for this result is not available because of this site's robots.txt  
Learn more

**MATLAB Primer, 8e - MATLAB & Simulink - MathWorks**  
<https://www.mathworks.com/support/books/book52873.html> ▾  
MATLAB Primer, covering the significant new features of MATLAB. This pocket book serves as a resource for students and engineers requiring a high-level ...

**Matlab Primer - UCSD Mathematics**  
[math.ucsd.edu/~driver/21d-s99/matlab-primer.html](http://math.ucsd.edu/~driver/21d-s99/matlab-primer.html) ▾  
by K Sigmon - Cited by 143 - Related articles  
MATLAB Primer Second Edition. Kermit Sigmon Department of Mathematics University of Florida.  
Department of Mathematics = University of Florida ...

**[PDF] MATLAB Primer - Department of Mathematics**  
[www.math.toronto.edu/mpugh/primer.pdf](http://www.math.toronto.edu/mpugh/primer.pdf) ▾  
The Third Edition of the MATLAB Primer is based on version 4.0/4.1 of MATLAB. ... serv@math.ufl.edu containing the single line send matlab/primer.tex.

**[PDF] MATLAB Primer - Yale Image Processing and Analysis Group**  
[ipagwww.med.yale.edu/~stalib/beng445/primer.pdf](http://ipagwww.med.yale.edu/~stalib/beng445/primer.pdf) ▾  
by K Sigmon - 1989 - Cited by 24 - Related articles  
The MATLAB Primer may be distributed as desired subject to the following con- ... MATLAB is an interactive, matrix-based system for scientific and engineering ...

See results about

**MATLAB Primer (Book by Kermit Sigmon)**  
Originally published: 1993  
Author: Kermit Sigmon

**MATLAB Primer, Eighth Edition (Book by Timothy A. Davis)**  
Originally published: December 29, 2004  
Author: Timothy A. Davis



# Math128A algorithms vs. General Purpose algorithms

For any vector  $\mathbf{x} \in \mathbb{R}^n$ , compute its norm

$$\|\mathbf{x}\|_2 = (x_1^2 + x_2^2 + \cdots + x_n^2)^{\frac{1}{2}} = \left( \sum_{k=1}^n x_k^2 \right)^{\frac{1}{2}}.$$

- ▶ **INPUT:**  $n, x_1, \dots, x_n$ .
- ▶ **OUTPUT:** Norm.
- ▶ **Step 1:** Set **SUM** = 0.
- ▶ **Step 2:** For  $k = 1, \dots, n$  set **SUM** = **SUM** +  $x_k * x_k$ .
- ▶ **Step 3:** Set **Norm** =  $\sqrt{\text{SUM}}$ .
- ▶ **Step 4:** Output **Norm**.  
STOP.

## Math128A algorithms vs. General Purpose algorithms

```
>> n=10;
>> x = (1:n)';
>> sum = 0;
>> for k=1:n
sum = sum + x(k)*x(k);
end
>> nrm = sqrt(sum);
>> disp([nrm, abs(nrm-sqrt(n*(n+1)*(2*n+1)/6))])
1.9621e+01  0.0000e+00
```

## Math128A algorithms vs. General Purpose algorithms

```
.... ...
>> n=10;
>> x = 1e200*(1:n)';
>> sum = 0;
>> for k=1:n
    sum = sum + x(k)*x(k);
end
>> norm = sqrt(sum);
>> disp([norm, abs(norm(x)-1e200*sqrt(n*(n+1)*(2*n+1)/6))])
    Inf  0.0000e+00
...
|
```